

Efficient RISC V Compute Platforms for Enabling the AI Revolution

Dr. Sujay Deb

CTO, AarFive Designs Pvt. Ltd.

Institute Chair Professor, Department of ECE & CSE

Indraprastha Institute of Information Technology Delhi (IIIT Delhi)



INDRAPRASTHA INSTITUTE of
INFORMATION TECHNOLOGY
DELHI



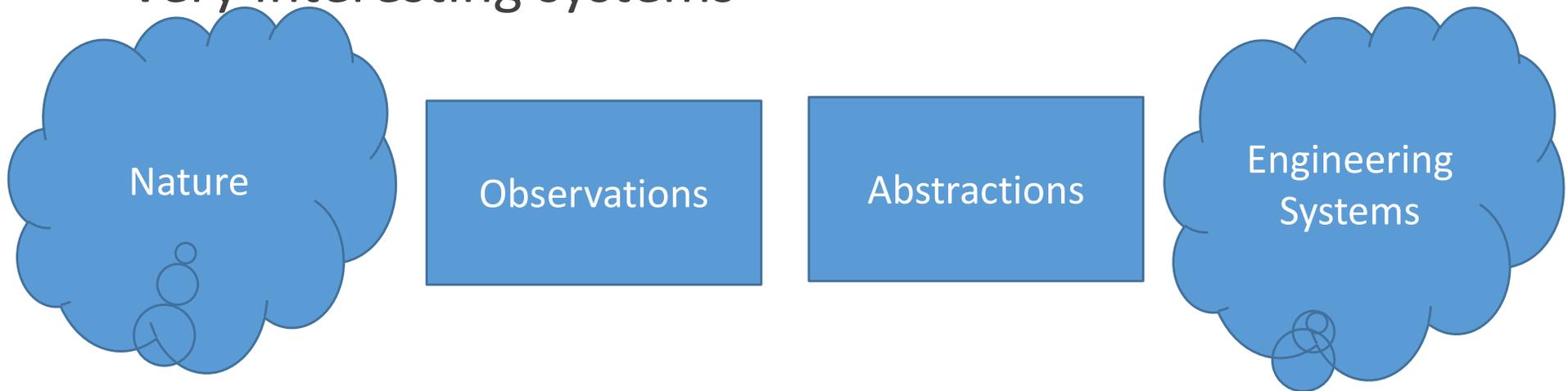
**AARFIVE
DESIGNS**



Engineering Approach



- What is Engineering?
 - the purposeful use of science
- We are here to employ the facts of nature to build very interesting systems



- Take complicated things, build layers of abstraction, and simplify things so that we can build useful systems.

Abstractions in Modern Computing Systems

Application

Application Requirements:

- Suggest how to improve architecture
- Provide revenue to fund development

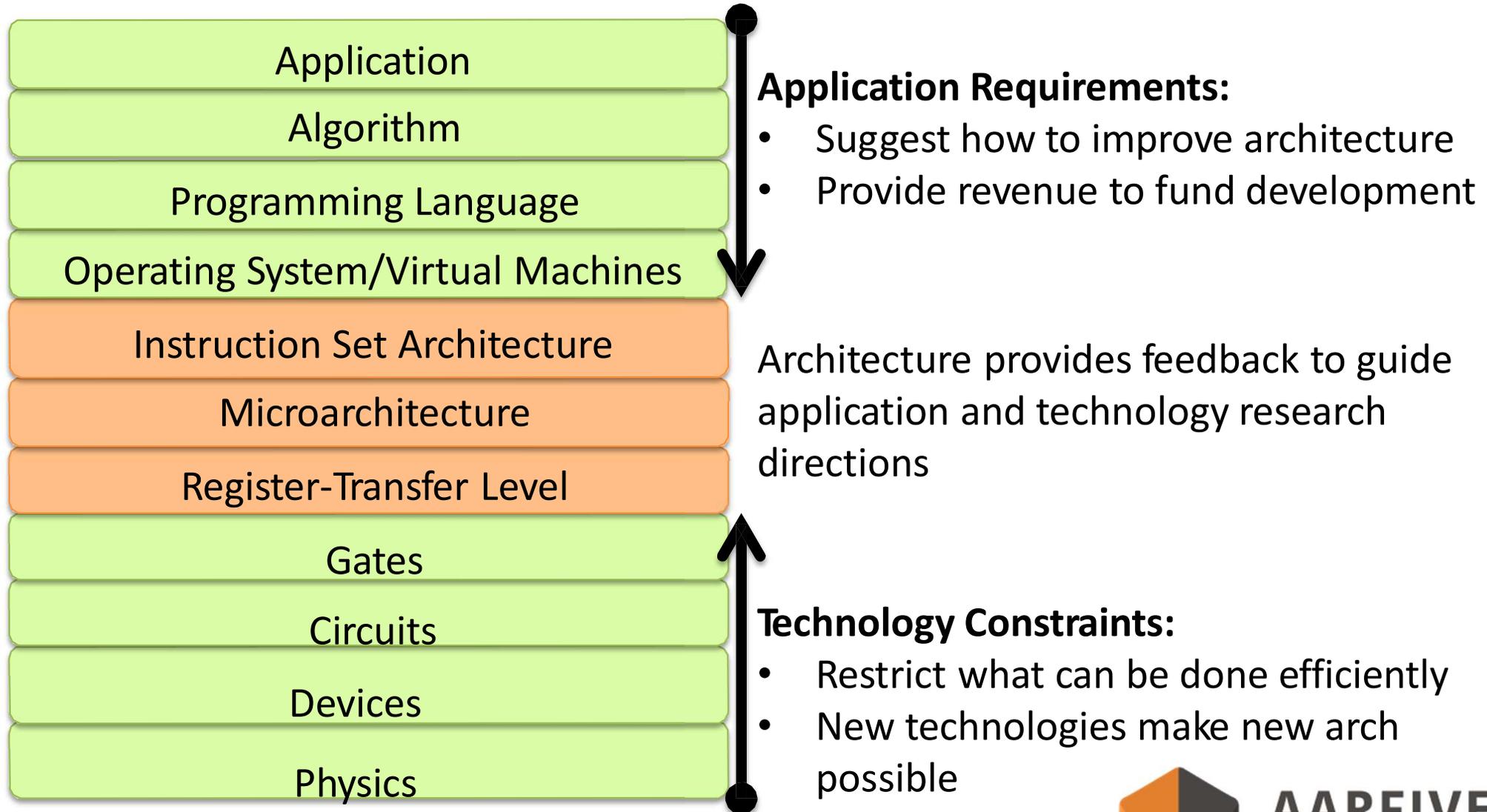
Technology Constraints:

- Restrict what can be done efficiently
- New technologies make new arch possible



**AARFIVE
DESIGNS**

Abstractions in Modern Computing Systems



Apple's M3 Max Chip Enters Testing Phase: 16-Core CPU, 40-Core GPU Onboard

Meta is developing its own AI chip - here's hoping it goes better than the Metaverse

News By Muskaan Saxena published May 19, 2023

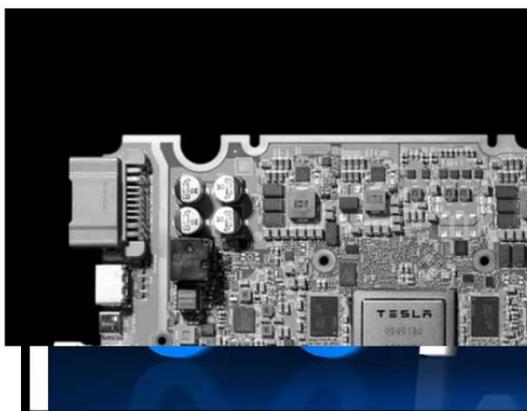
M2 Max. The chip features a models. Hints from test logs enter process.



Tesla's new self-driving chip is your best

Google's Future Tensor Chips, Fully In-House Design, Independent of Samsung

By Sofie Gallardo - July 10, 2023 906 0



Apple's M3 Max chip is expected to power the 17-inch MacBook Pro models in 2024.



**AARFIVE
DESIGNS**

NVIDIA and MediaTek's new Arm-based AI chip rumored for 1H 2025 to fight Intel, AMD, Qualcomm

6-Core CPU, 40-

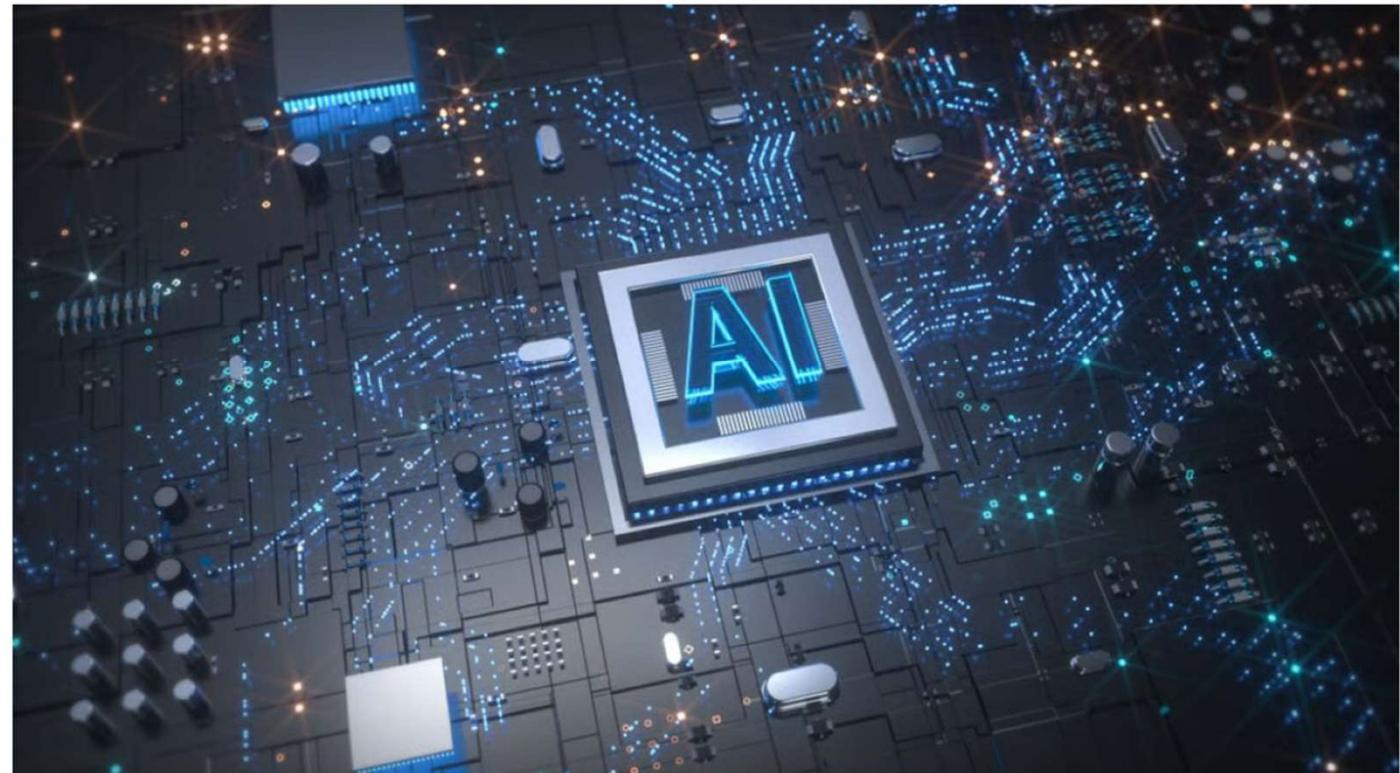
MediaTek and NVIDIA expected to unveil their Arm-based AI processor in the first half of 2025, will

Google and Meta unveil their Nvidia challengers

Custom AI chips can give Nvidia headaches

By Mayank Sharma

April 23, 2024



Meta and Google have just shown the way to take on Nvidia in the AI chips race

Meta and Google just announced their intention to push further into the AI chip race, currently dominated by Nvidia. Both companies are touting their new custom AI chips as a key element of their AI strategies.

Te
th



Publ

Apple's M3 Max chip is expected to power the 17- and 19-inch MacBook Pro models in 2024

NEWS SEMICONDUCTORS

Expect a Wave of Wafer

allows for one version no
in 2027

'Feels like magic!': Groq's ultrafast LPU could well be the first LLM-native processor — and its latest demo may well convince Nvidia and AMD to get out their checkbooks

News

By Wayne Williams published 28 February 2024

NEWS FEED

Google's Tensor G5 processor to enter tape-out stage, manufactured with TSMC's 3nm process

by TechNode Feed Jul 2, 2024

onal bottlenecks

ate commission. [Here's how it](#)

News

NVIDIA plans ARM chip for laptops: A potential Apple M4 killer?



Philipp Briel · 13. August 2024



Popular Posts

- > Backforce V test: New refer in the 400 euro segment
- > Backforce One Plus review: the competition look old aga
- > Backforce One (Plus) Discou Coupon Codes
- > Backforce One review: Can convince with a gaming cha
- > Backforce versus noblecha thrones for a hallelujah

AARFIVE DESIGNS

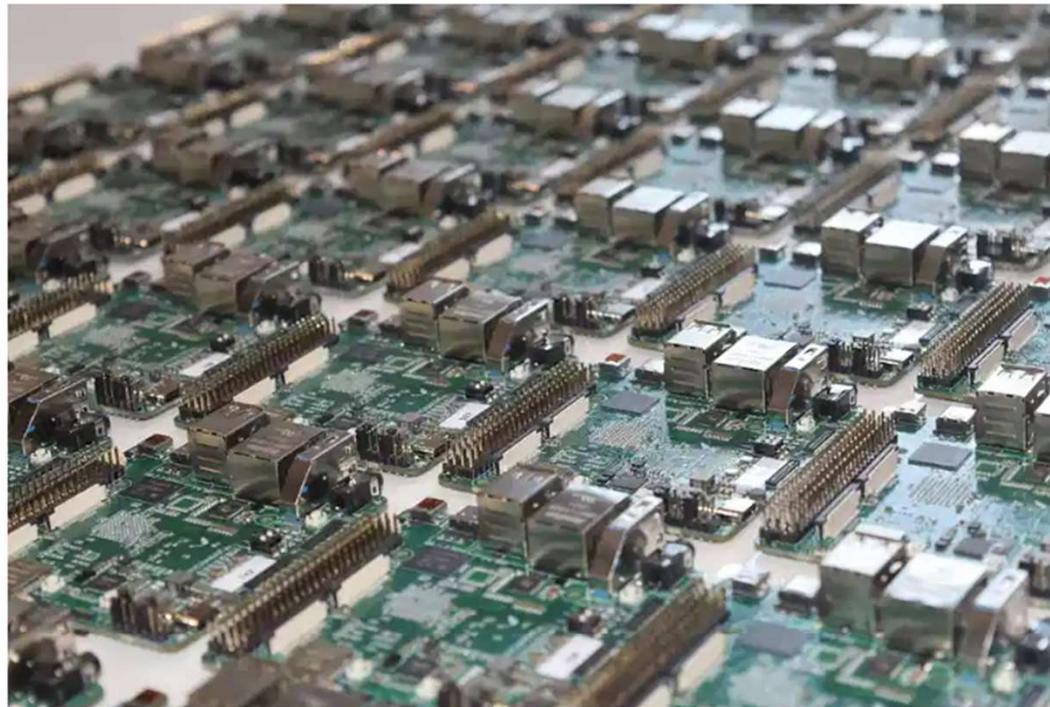
TSMC's wafer-scale integration tech is the key to TESLA

- > Backforce versus Can techa training chairs (almost) at e
- > Backforce versus Secret al versus e-ports experience
- > Coupons & Discount Codes

Innovation At The 'Edge': Hardware Startup From Surat Pioneers The Reconfigurable Edge Computer — And Reshapes Computing As We Know It

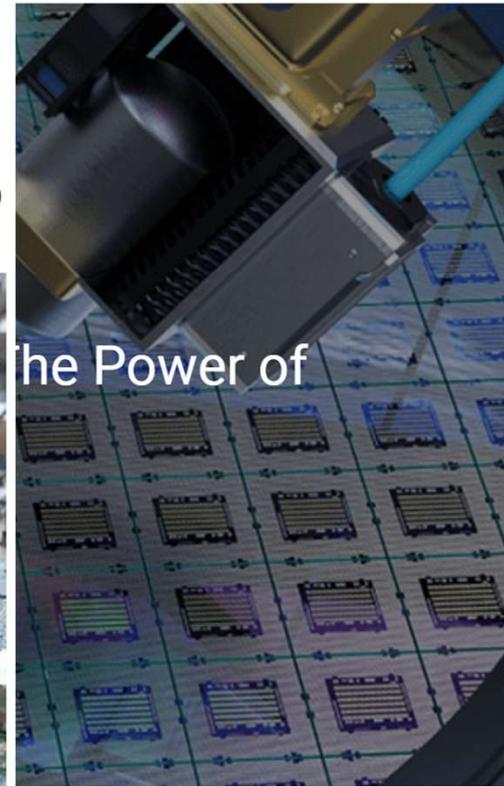
KARAN KAMBLE

May 25, 2024, 02:02 PM | Updated 02:02 PM IST



Vicharak's Vaaman edge computing board

Last month, Surat-based hardware startup Vicharak launched their very first product — an edge computing board called Vaaman.



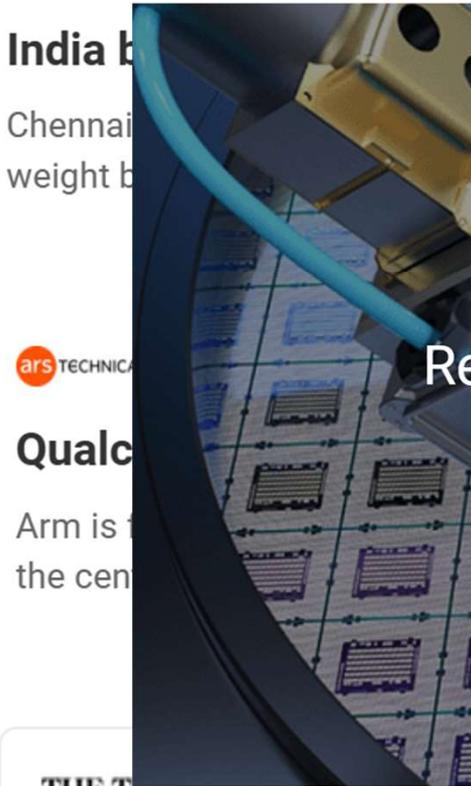
The Power of

eled by strategic or (DIR-V) Programme. conductor design, instead of



**AARFIVE
DESIGNS**

Indiatim **TESSOLVE**
A HERO ELECTRONIX VENTURE



India b
Chennai
weight b

ars TECHNICAL

Qualc

Arm is
the cen

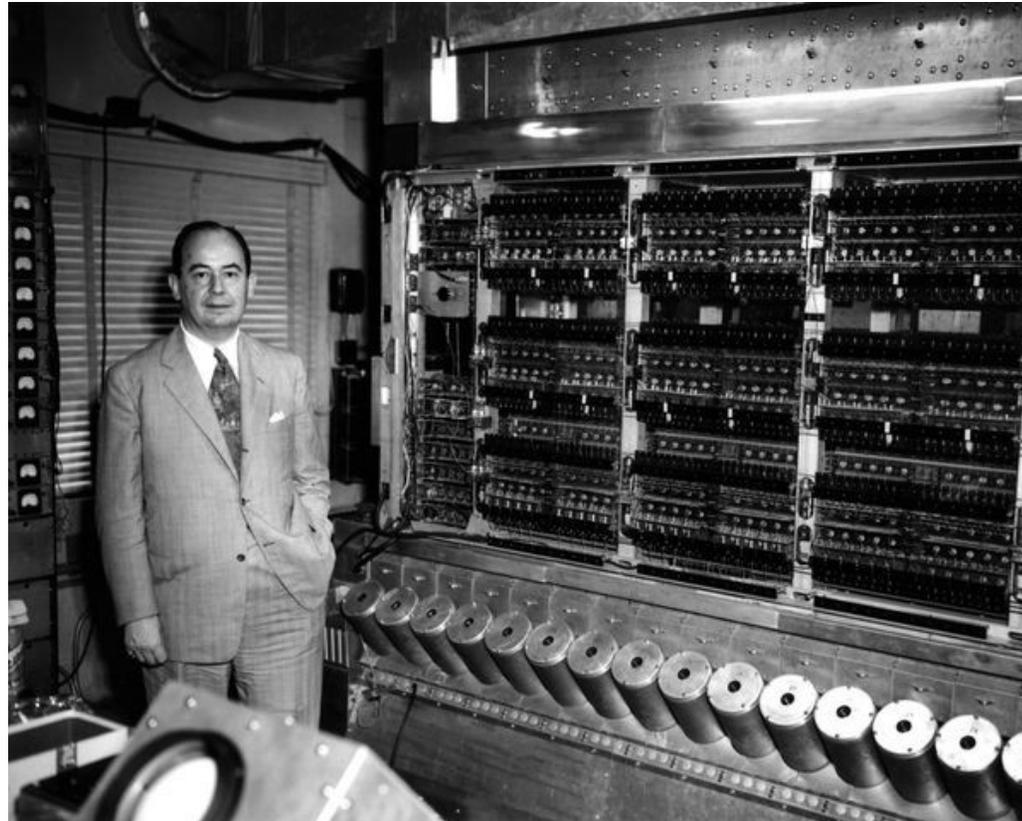
THE T

Expl

India is promoting the use o
ARM technology that is ...

Indic
part

Computers Then..



John von Neumann with the IAS Computer (Courtesy of the Shelby White and Leon Levy Archives Center, Institute for Advanced Study (IAS))

Computers Now



**AARFIVE
DESIGNS**

Moore's Law



VISUALIZING PROGRESS

If transistors were people

If the transistors in a microprocessor were represented by people, the following timeline gives an idea of the pace of Moore's Law.



Now imagine that those 1.3 billion people could fit onstage in the original music hall. That's the scale of Moore's Law.

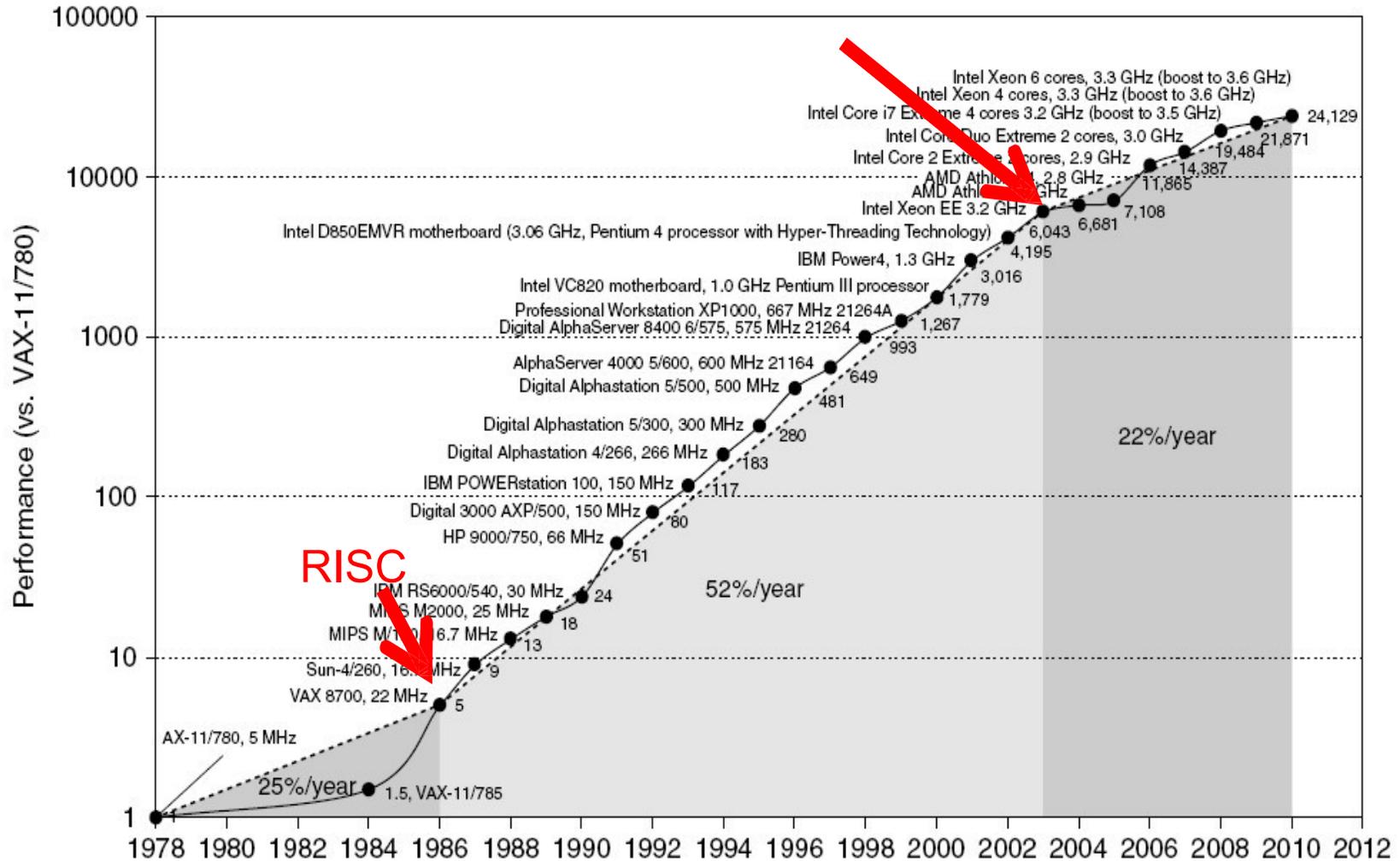
Courtesy:

<http://www.intel.com/content/www/us/en/silicon-innovations/moores-law-technology.html>

Sequential Processor Performance



Move to multi-processor



From Hennessy and Patterson Ed. 5 Image Copyright © 2011, Elsevier Inc. All rights Reserved.



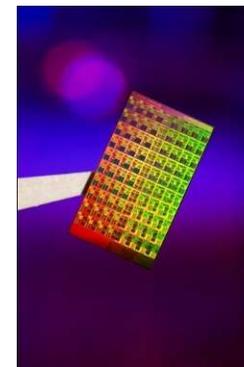
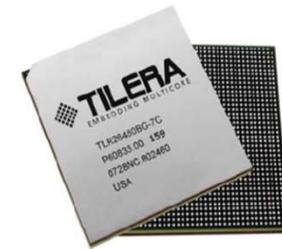
The era of Many-Core systems



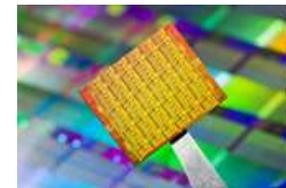
- How to keep up with demands on computational power?
 - Can not scale clock frequency
 - Solution: Increase number of cores - parallelism
 - Mass Market production of Intel, AMD dual-core and quad-core CPUs
 - Custom Systems-on-Chip (SoCs)
 - Many Core chips from Tiler for networking, cloud computing and multimedia applications.



Adapteva's
Epiphany



Intel 80 core
processor



Single-chip
Cloud
Computer

'Number of cores will double every 18 months'

- Prof. A. Agarwal, MIT, founder of Tiler Corporation



**AARFIVE
DESIGNS**

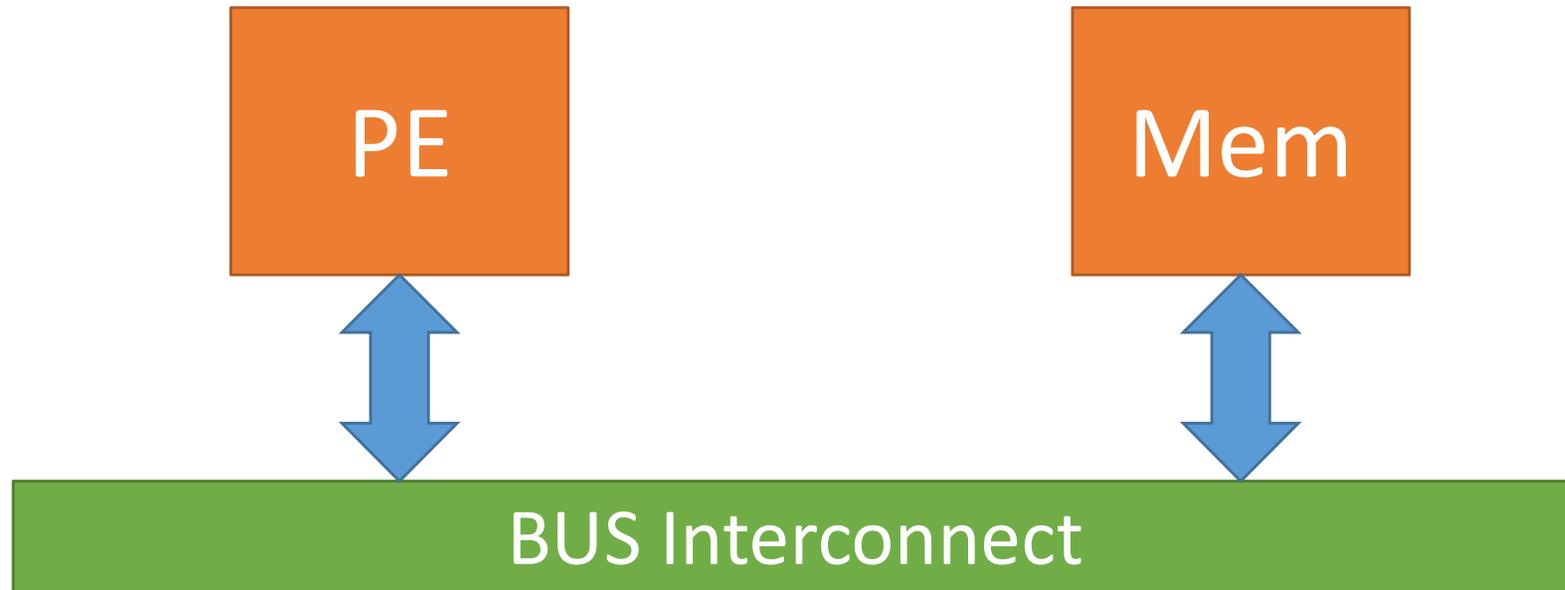
The era of Many-Core systems



- We are at the early stage of Many-core Processor evolution
 - Many-core is going to be ubiquitous
- Immense possibilities:
 - Server-type performance on handheld devices

	
ASCI Red: 1 TF	Knights Corner: 1 TF
1997 First System 1 TF Sustained	2011 First Chip 1 TF Sustained
9298 Pentium II Xeon	1 22nm Chip
OS: Cougar	OS: Linux
72 Cabinets	1 PCI express slot

Basic Computing System

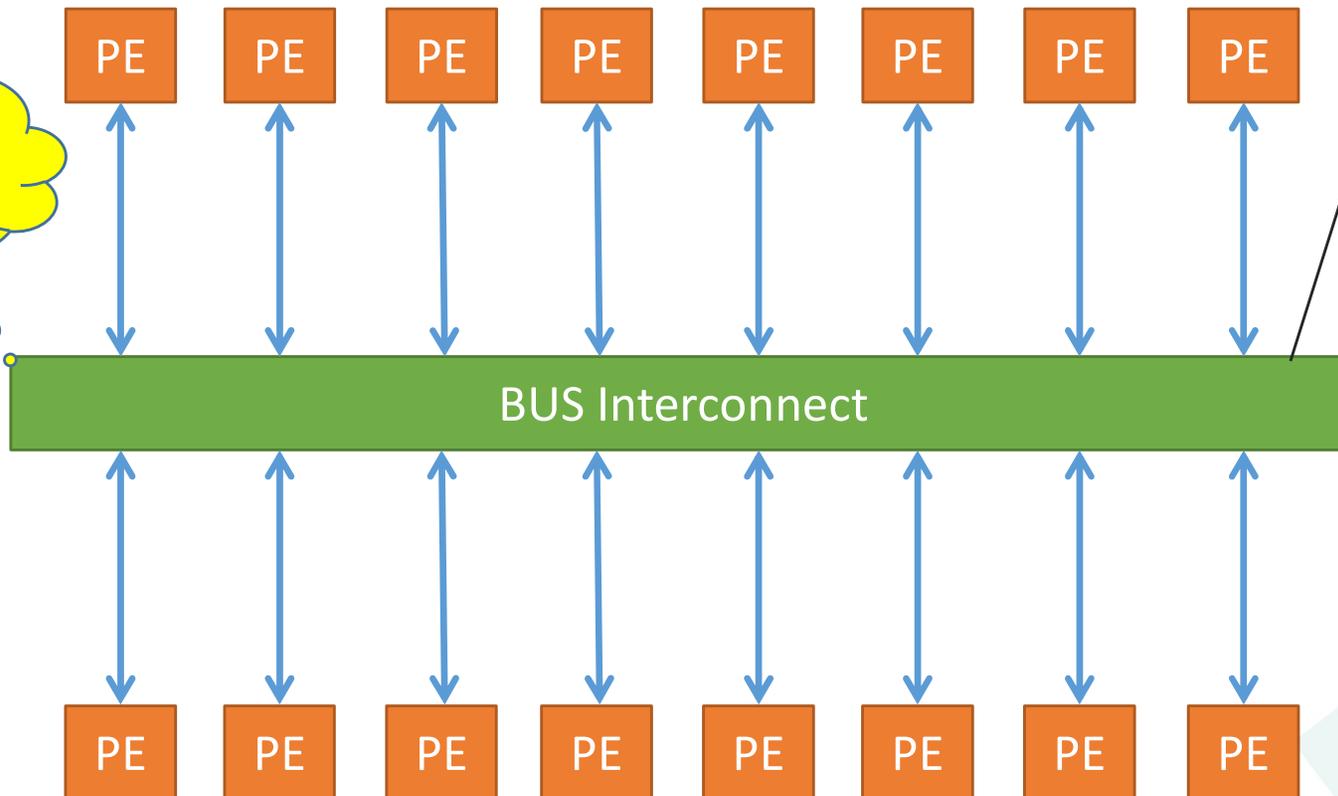


Multi-Core Systems

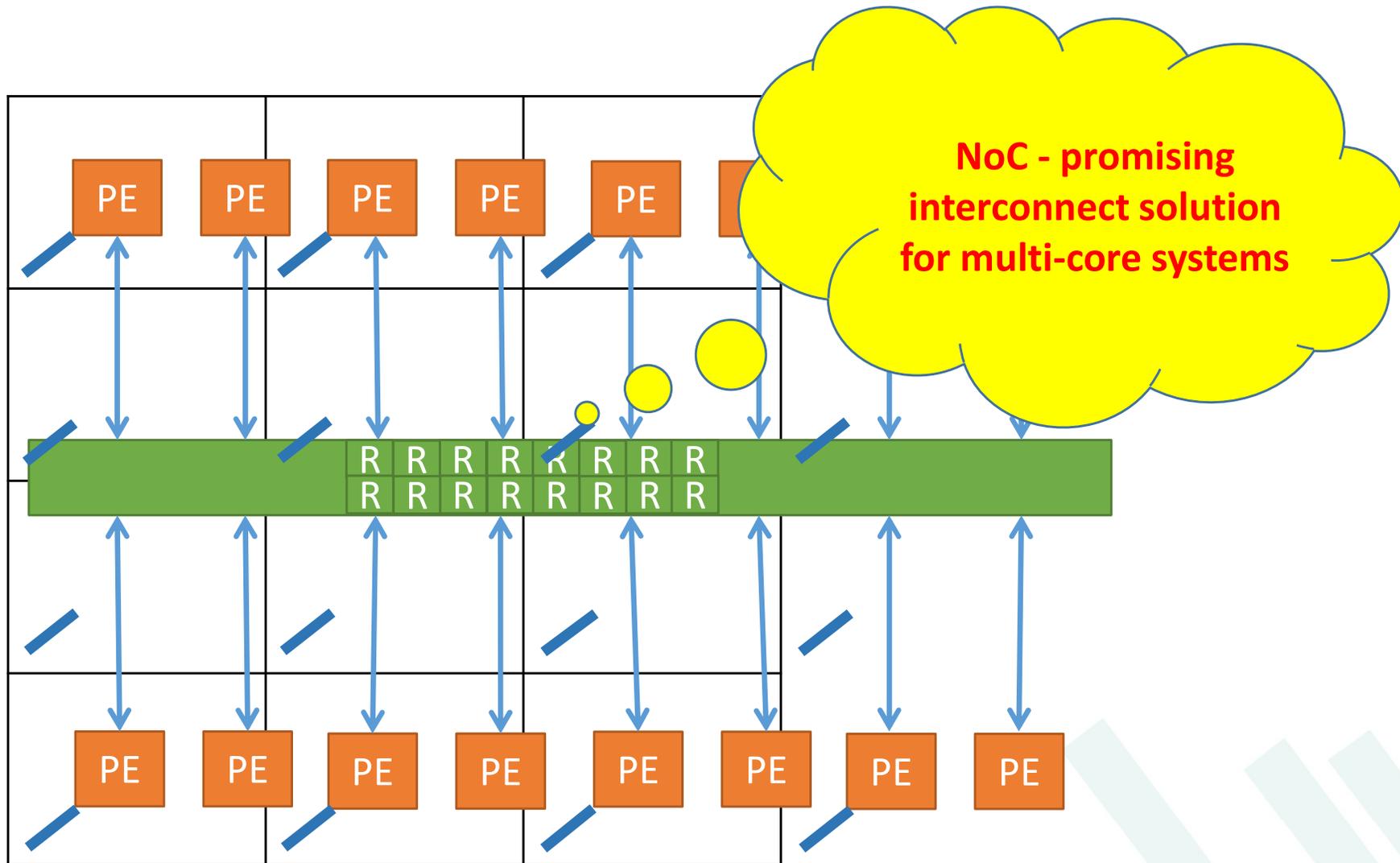


Not Scalable

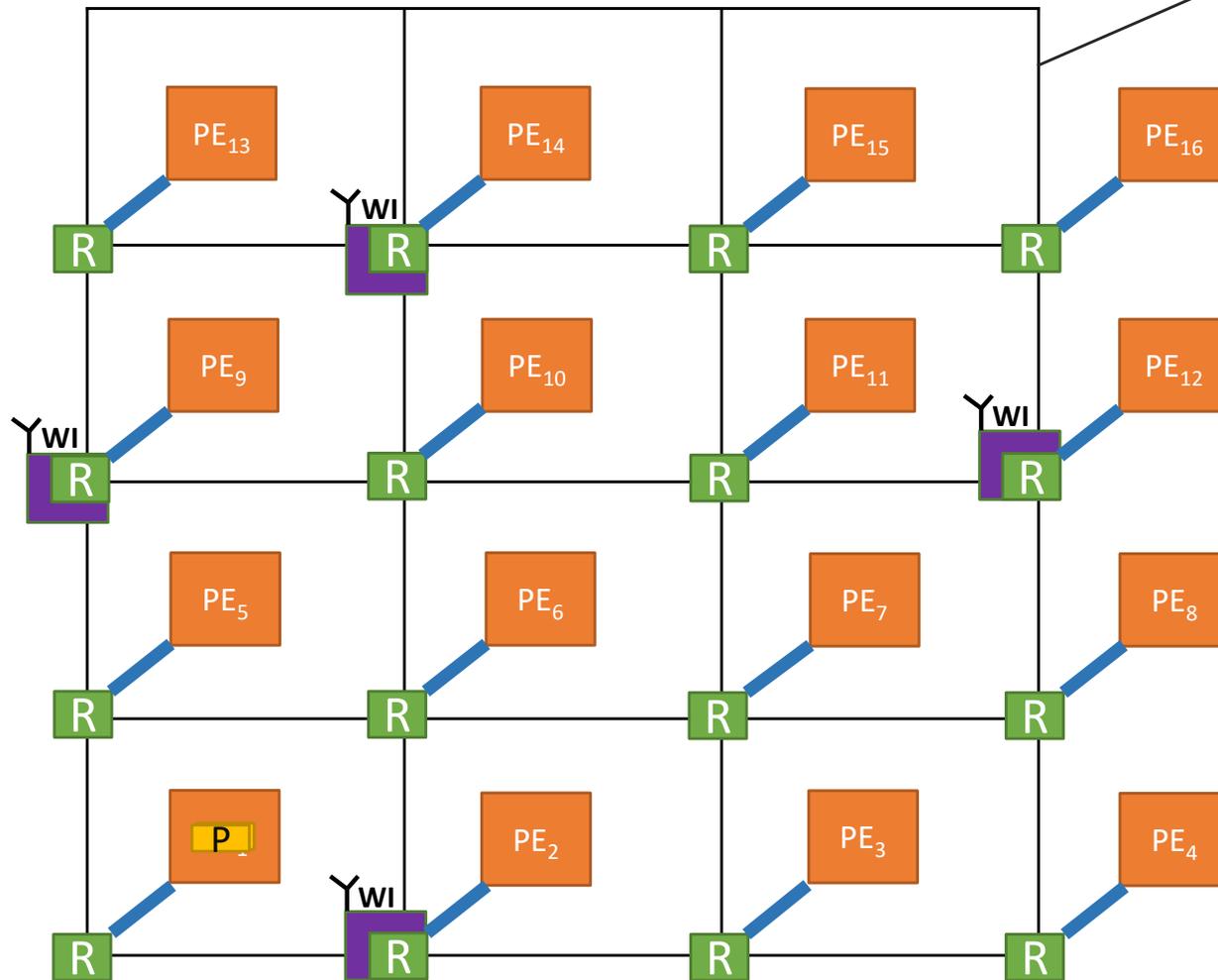
Kilocore



Network-on-Chip (NoC)

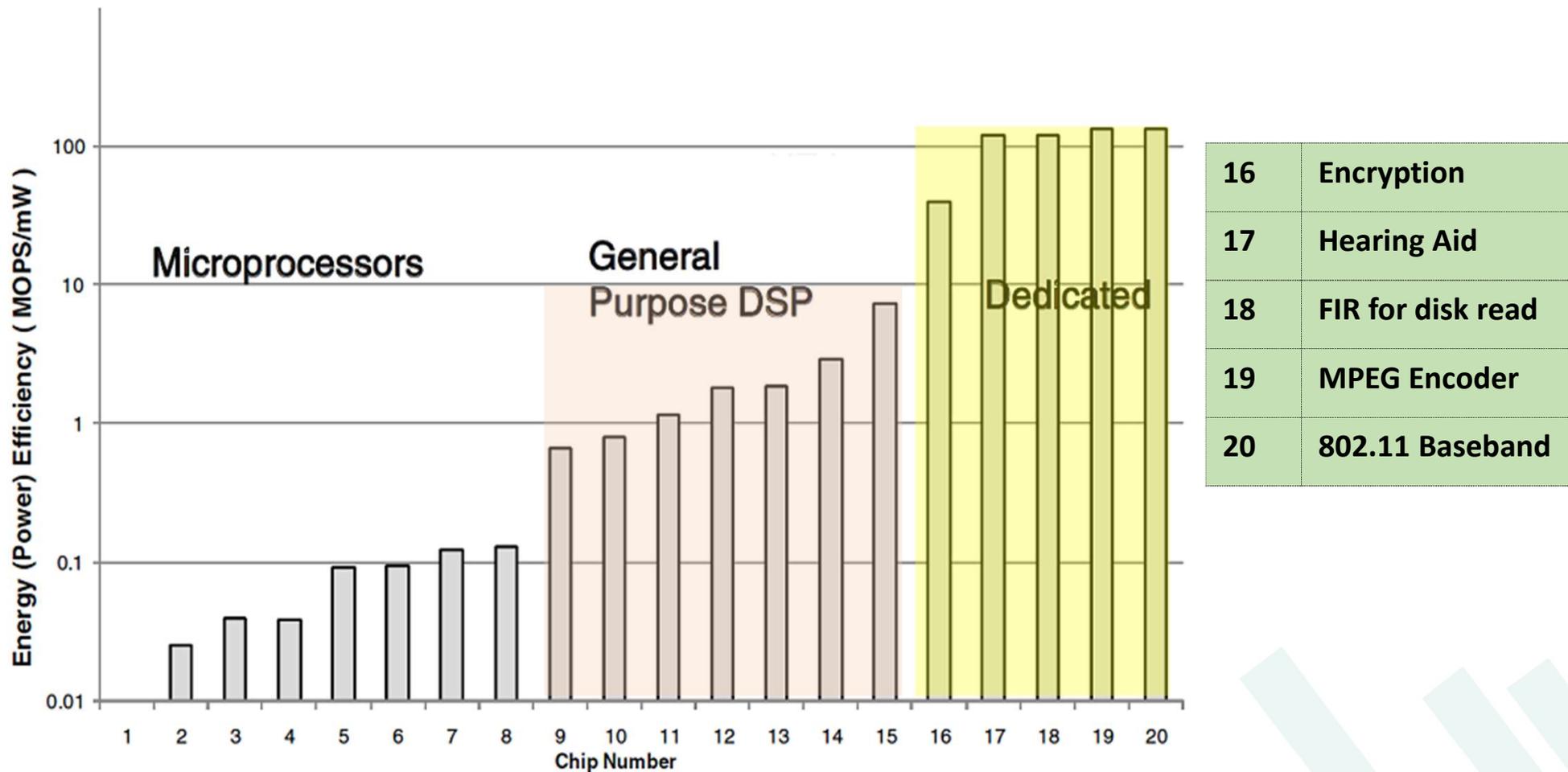


Wireless NoC (WNoC)



- Multi hop communication
- Power and performance bottleneck

Efficient System Design Strategy



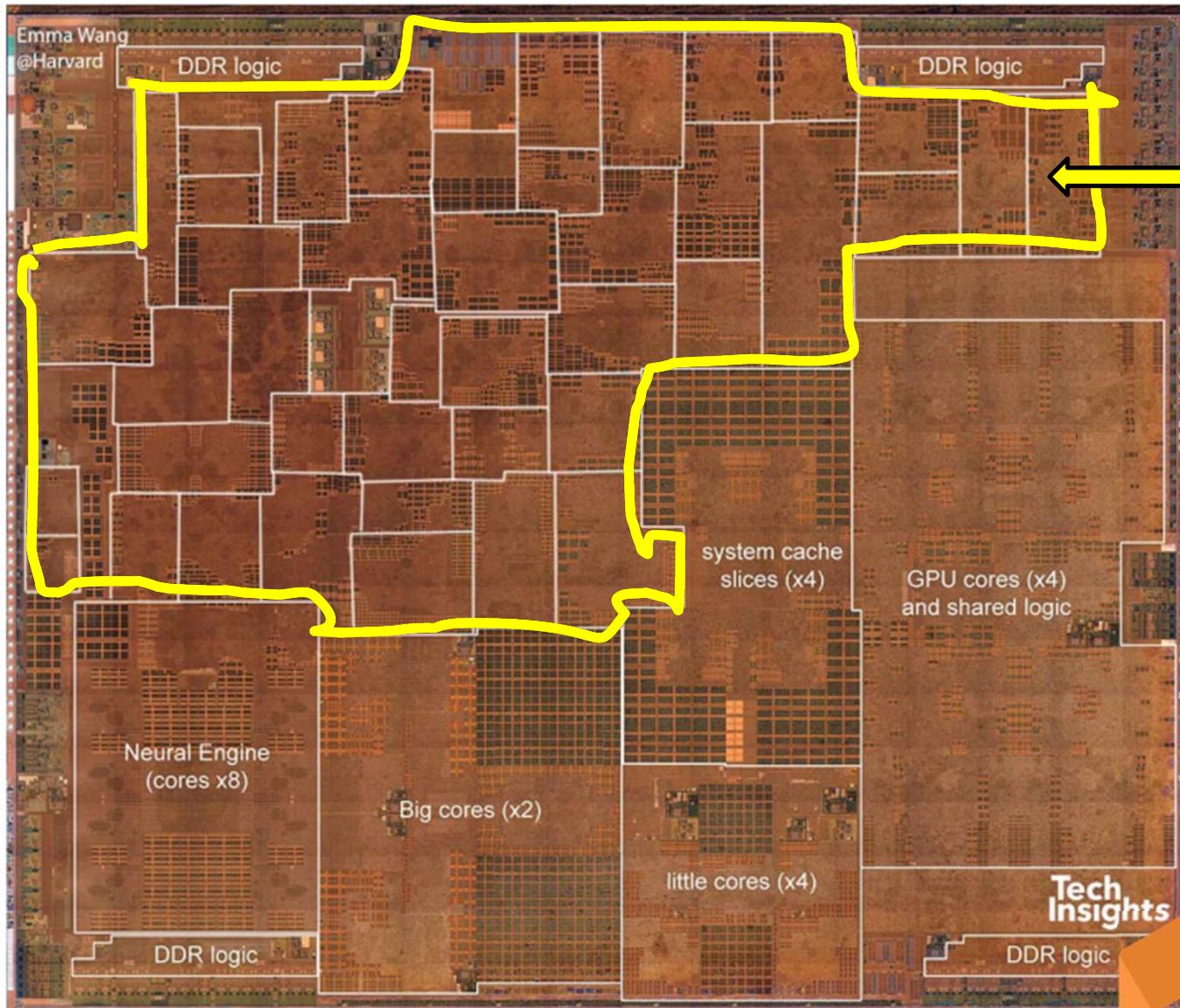
16	Encryption
17	Hearing Aid
18	FIR for disk read
19	MPEG Encoder
20	802.11 Baseband

Fig. 1. Efficiency of specialized accelerators

Efficient System Design Strategy



Apple System-on-Chip : A12



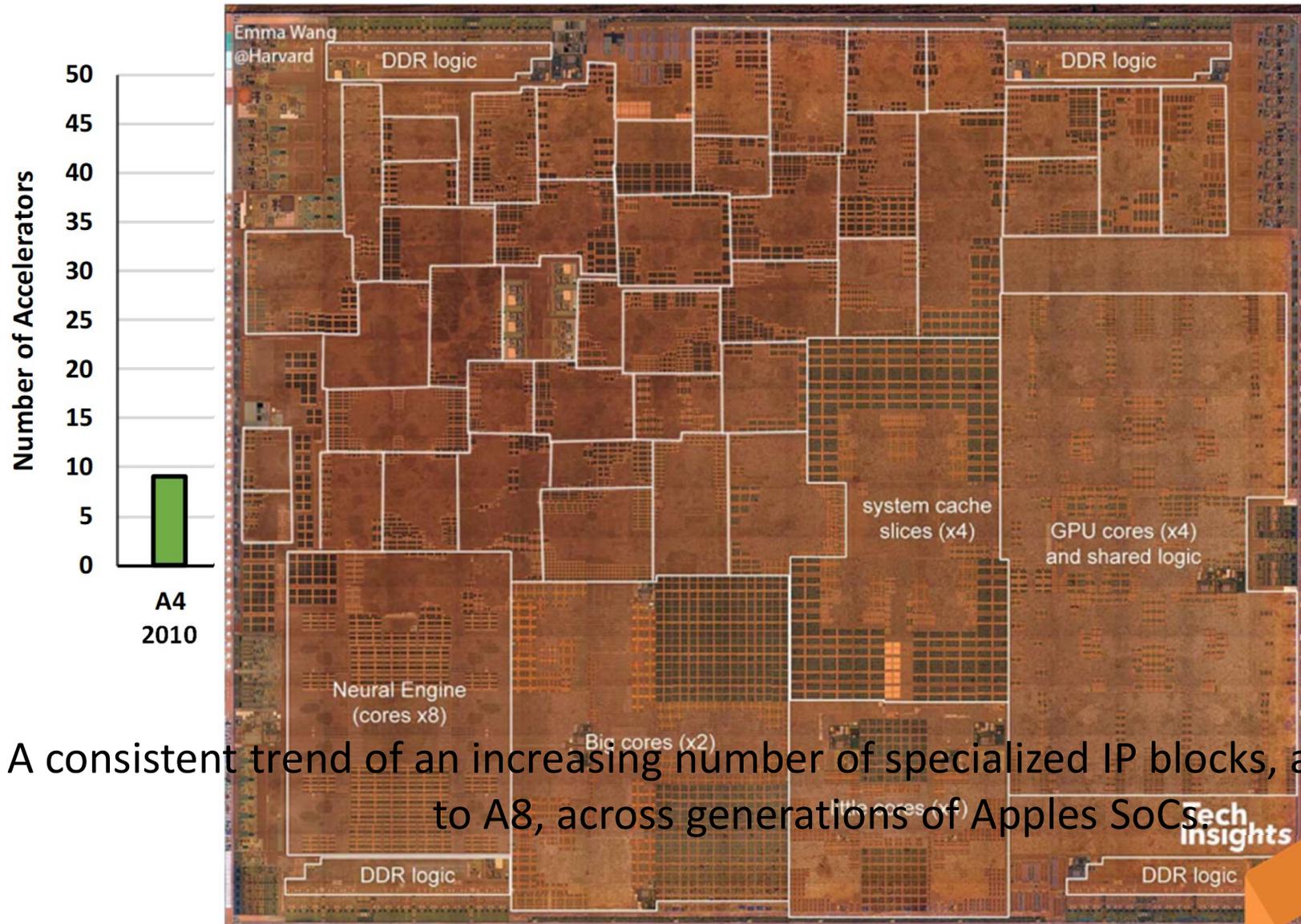
September 16, 2024

[Pic source: <http://vlsiarch.eecs.harvard.edu/research/accelerators/die-photo-analysis/>]



**AARFIVE
DESIGNS**

Efficient System Design Strategy



A consistent trend of an increasing number of specialized IP blocks, almost 3x from A4 to A8, across generations of Apples SoCs

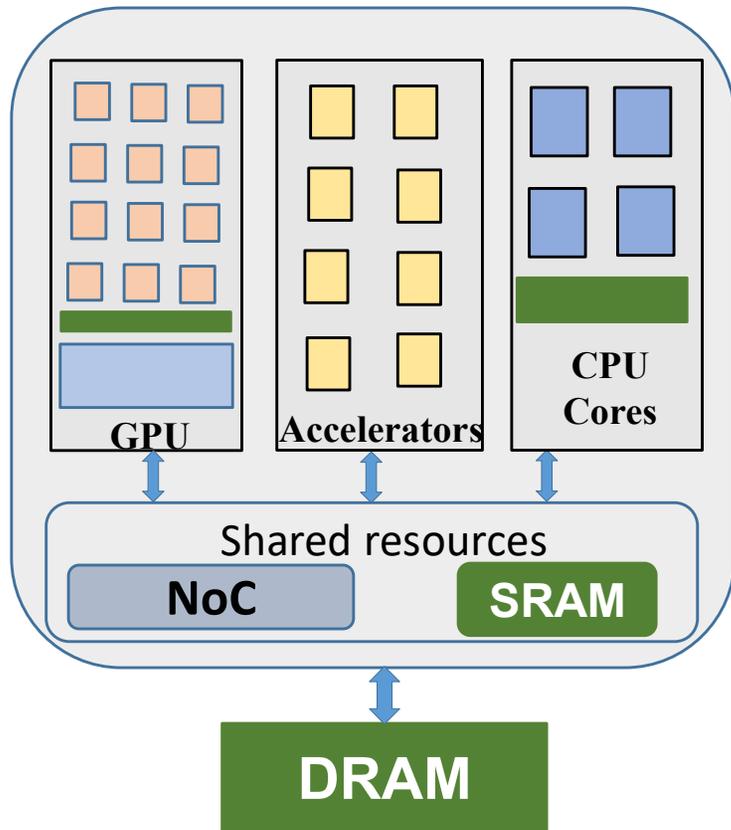
September 16, 2024

[Pic and data source: <http://vlsiarch.eecs.harvard.edu/research/accelerators/die-photo-analysis/>]

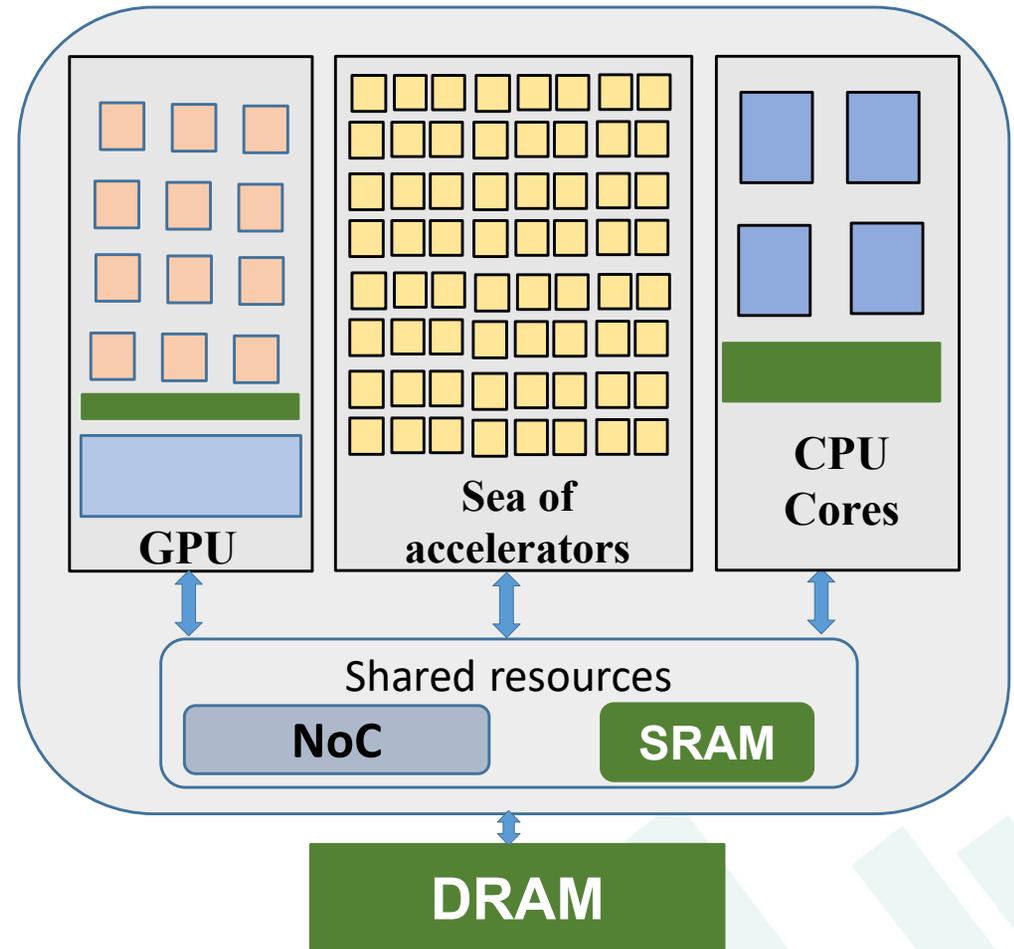


**AARFIVE
DESIGNS**

Future of SoC

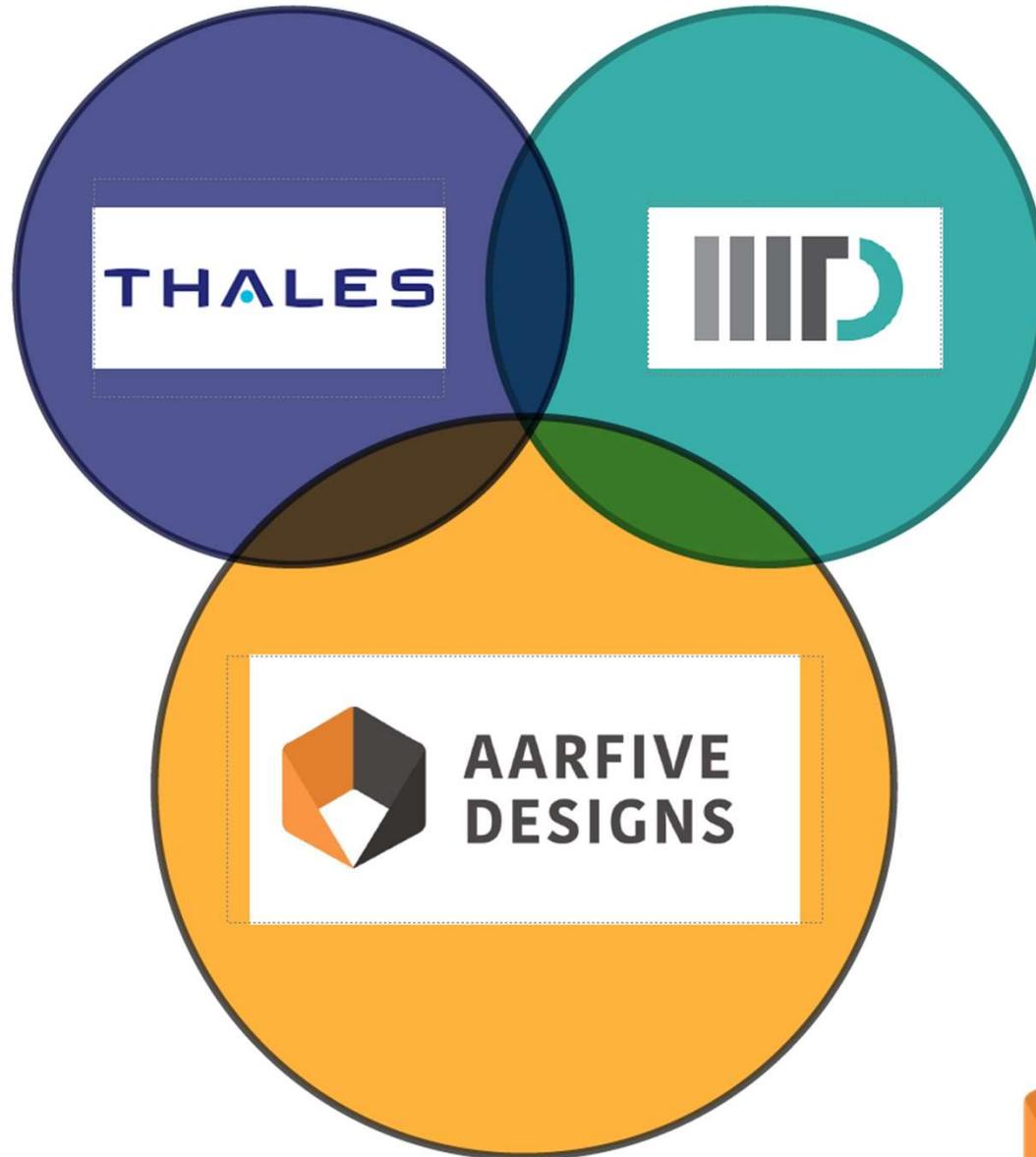


Present accelerator-rich SoC

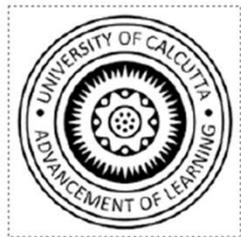


Emerging accelerator-rich SoC

About AarFIVE Design Pvt Ltd



People who **trust** us



**Academy
Partners**



**IP
Partner**



Thank You!



AARFIVE DESIGNS



INDRAPRASTHA INSTITUTE of
INFORMATION TECHNOLOGY
DELHI

Contact Us

 Website
www.aarfive.in

 Email
business@aarfive.in

 Address
IIIT - Delhi, near Govind Puri Metro
Station, Shyam Nagar, Okhla Industrial
Estate, New Delhi, Delhi, India

September 17, 2024